

Stat 602 Project Report

Kiegan Rice and Nate Garton

Due ???

Introduction

The goal of this competition was to try to train an effective binary classifier on a small training set with many predictor variables. The criteria used to determine the best solution was the area under the ROC curve generated with predictions on the test set. The training data set had dimensions 250×300 , and the test data had dimensions 19750×300 .

Exploratory Data Analysis

Data Quality

The quality of the data was pristine. There were no missing values in the test or training data. There was some class imbalance in the training set (roughly 64% of the observations were class 1). A 95%, two-sided confidence interval based on large sample normality assumptions for the population class proportion was $(0.58, 0.70)$. Thus, assuming that we were truly given a random sample from the total data, it seemed unlikely that the classes were truly balanced. However, even at the upper bound of the confidence interval, the class imbalance would not likely be large enough to suggest problems with typical classification models/algorithms.

There were 300 predictor columns given; however, there was no contextual information given about the predictors or the response. There was no information on whether columns represented any particular measured quantity, and each column was simply named $'0', '1', '2', \dots, '299'$. This makes exploratory data analysis particularly difficult, especially since significant dimension reduction needs to occur to avoid overfitting to noise in the data, and since the training data set only has 250 cases. One clear approach in this scenario is to investigate distributions of each predictor variable, particularly class-conditional distributions since we are working with a classification problem.

Class Conditional Distributions

Because there are so many predictor variables, it is hard to pull useful information from tables of numerical summaries. Figure 1 presents histograms of four summary statistics: minimums, maximums, means, and standard deviations for each column conditional on the class. We can see that the ranges of all of the columns are similar between and within each class. We can also see that columns means for both classes are centered at zero and look fairly Gaussian. There may be one or two “outlier” means in class 0 that are unusually small, but given the number of columns in the data, it was not immediately clear to us that this wasn’t simply an

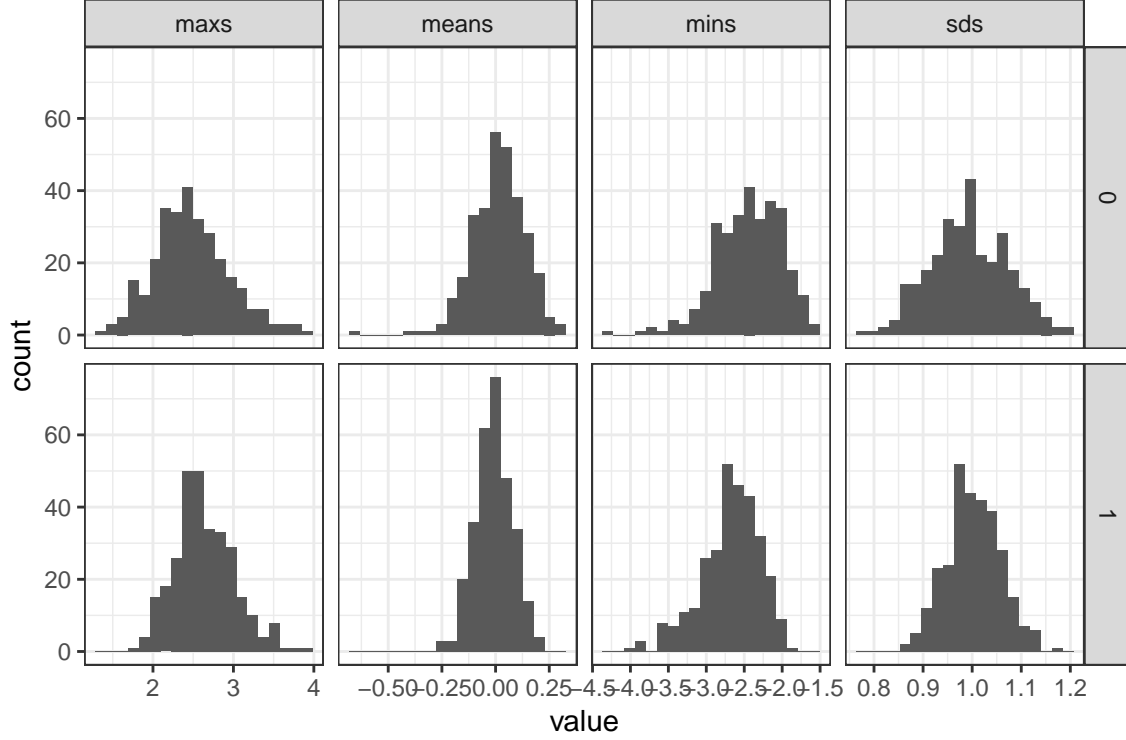


Figure 1: Class conditional histograms for column minimums, maximums, means, and standard deviations. Each row contains the column distribution for these four summary statistics for one of the two classes.

artifact of randomly sampled data. Standard deviations for both classes are centered at 1, and most standard deviations for both classes are between 0.9 and 1.1.

The behavior of these summary statistics seems consistent with data where most columns are $N(0, 1)$. With this in mind, we performed the Shapiro-Wilk test for normality on each column for each class. Assuming that the columns are independent for both classes, the distribution of p-values for each class should be approximately uniform. To test this, we performed Kolmogorov-Smirnov (KS) tests for the p-values from the Shapiro-Wilk tests. The two p-values from the (KS) tests were ≈ 0.904 and ≈ 0.997 for classes 0 and 1, respectively. This led us to believe that assuming the normality of the predictor variables was likely reasonable.

If we could also conclude that the predictor variables within each class were uncorrelated, then this would allow us to effectively model the class conditional densities. To assess this, we first looked at histograms of the correlations. Figure 2 shows histograms of column correlations within each class. We see that both distributions are centered at essentially zero and appear to have a bell shape. Such a pattern is consistent with the sampling distribution of the correlation between two uncorrelated Gaussian random variables. The distribution of correlations for class 0 appears to be a bit more diffuse, but it is unclear exactly what to make of that. We performed z-tests after using the Fisher transformation for each correlation and then performed the KS test to test whether the p-value distributions for each class were close to uniform. The p-values from these tests were ≈ 0.113 and ≈ 0.905 for class 0 and class 1, respectively. We take these tests with a grain of salt because the correlations used in the test are not independent. However, if there were a handful of high correlations between several predictor variables, we might expect to see more p-values near 0. Finally, the

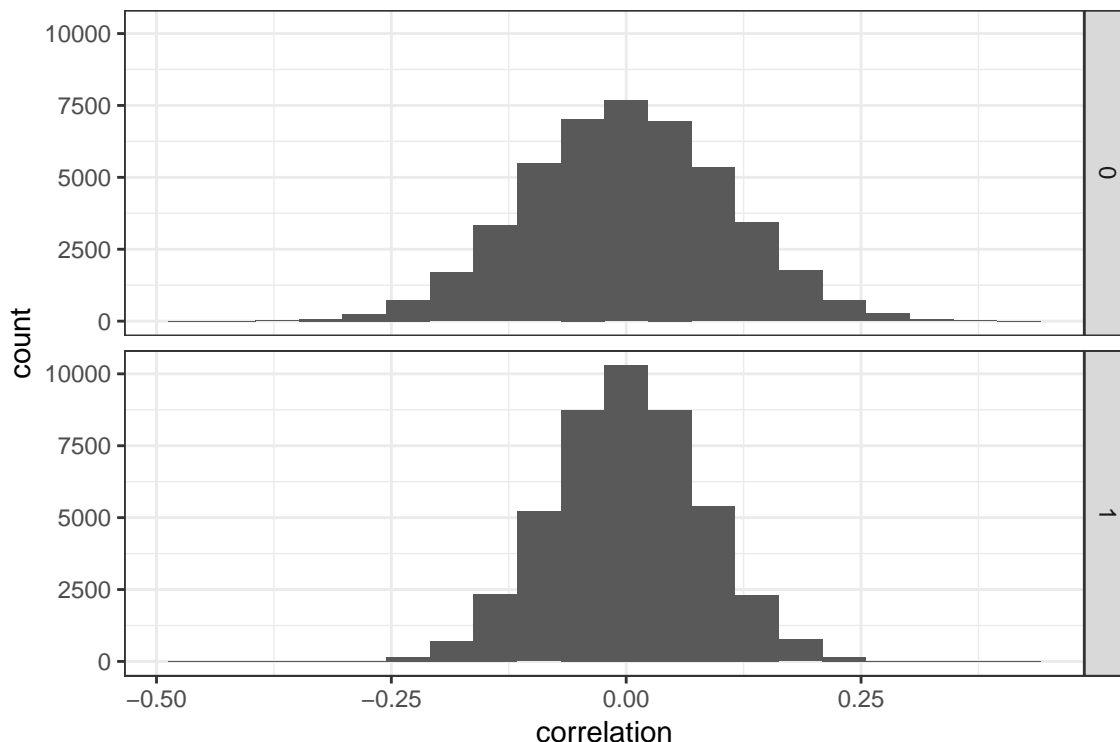


Figure 2: Histograms of column correlations for each class.

fact that the correlation distribution for class 0 seems more diffuse than class 1 suggests that there might be some subtle differences in correlation structure between the two classes. However, performing a two-sample KS test comparing the correlation distributions between the classes results in a p-value of ≈ 0.251 . Again, the assumption of independence is violated, so we take this with a grain of salt. But, this serves to reinforce the idea that any difference in correlation structure between the two classes is subtle, if it exists.

Visually investigating class-conditional densities for each column also gives a good look at columns where there may be some slight separation between the two response classes, seen in Figure 3. However, additional statistical tests should be done to assess whether the densities are significantly different, especially since we have such a large number of columns. Some visual separation could be an artifact of random sampling for the training set.

look at t-tests between column means

Feature Engineering

Creative feature engineering proved particularly difficult for this problem, especially due to the murky nature of the data. While data were pristine in the sense of completeness, the lack of contextual information about predictors makes context-based feature engineering impossible. The large number of predictors also poses the challenge of weeding out noise without removing useful information; dimension reduction plays the first big role in constructing a good set of features.

LASSO is a natural choice for dimension reduction here, as it will shrink many predictor parameter estimates

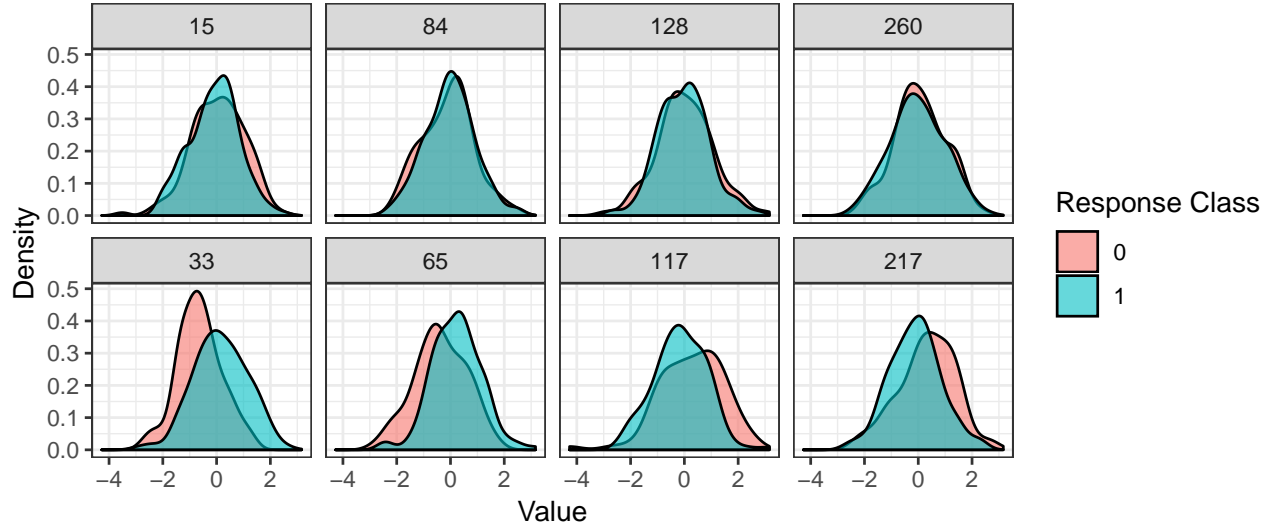


Figure 3: Example of class-conditional densities for several columns in the dataset. The first row shows class-conditional densities that do not demonstrate any separation, which is very common in the dataset. The bottom row demonstrates class-conditional densities observed in the training data which seem to show separation.

to zero in favor of those that more adequately describe the response variable. A logistic regression LASSO model using all 300 predictors was fit using the `glmnet` package, in particular, the `cv.glmnet` function with $\alpha = 1$.

Prediction Methods

Conclusions